



(11)

EP 0 935 210 A2

(12)

EUROPEAN PATENT APPLICATION

(43) Date of publication:
11.08.1999 Bulletin 1999/32

(51) Int Cl.⁶: **G06F 19/00**

(21) Application number: 99300900.0

(22) Date of filing: 08.02.1999

(84) Designated Contracting States:
AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU
MC NL PT SE
Designated Extension States:
AL LT LV MK RO SI

(72) Inventor: Mack, David H.
Menlo Park, C.A. 94025 (US)

(74) Representative:
O'Connell, David Christopher et al
Haseltine Lake & Co.,
Imperial House,
15-19 Kingsway
London WC2B 6UD (GB)

(30) Priority: 09.02.1998 US 20743

(71) Applicant: **Affymetrix, Inc.**
Santa Clara, CA 95051 (US)

(54) Computer aided visualisation of expression comparison

(57) Innovative systems and methods for visualizing information collected from analyzing samples are provided. The samples may include nucleic acids, proteins, or other polymers. Gene expression level as determined from analysis of a nucleic acid sample is one possible analysis result that may be visualized. In one embodiment,

ment, a computer system may display the expression levels of multiple genes simultaneously in a way that facilitates user identification of genes whose expression is significant to a characteristic such as disease or resistance to disease. Additionally, the computer system may facilitate display of further information about relevant genes once they are identified.

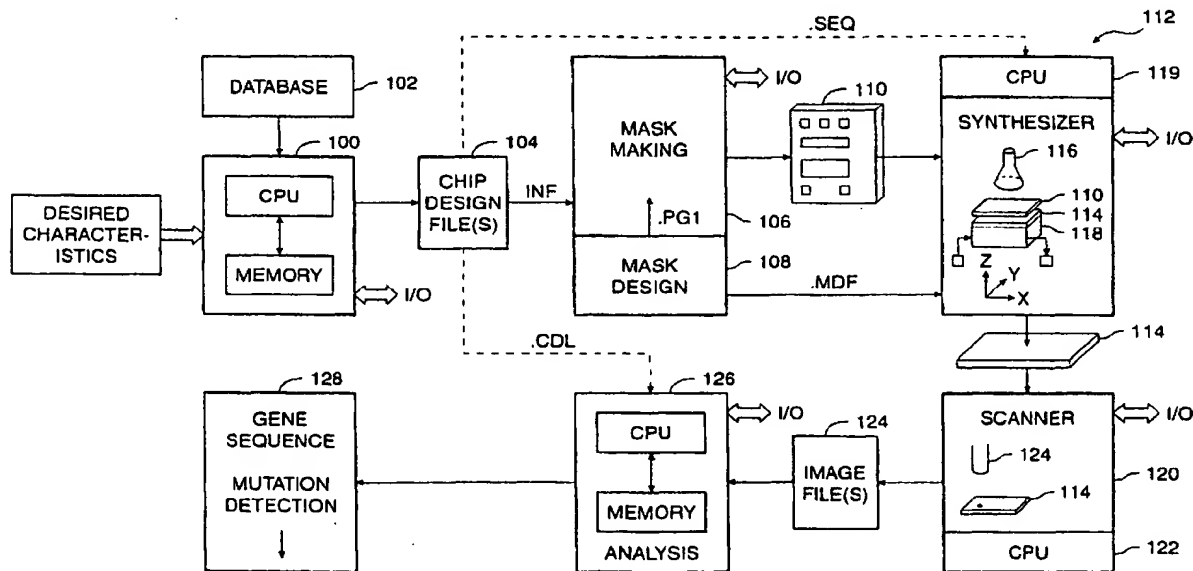


FIG. 3

Description

BACKGROUND OF THE INVENTION

[0001] The present invention relates to the field of computer systems. More specifically, the present invention relates to computer systems for visualizing analysis results.

[0002] Devices and computer systems for forming and using arrays of materials on a substrate are known. For example, PCT Publication No. WO 92/10588, incorporated herein by reference for all purposes, describes techniques for sequencing or sequence checking nucleic acids and other materials. Arrays for performing these operations may be formed according to the methods of, for example, the pioneering techniques disclosed in U. S. Patent No. 5,143,854 and U.S. Patent No. 5,593,839 both incorporated herein by reference for all purposes.

[0003] According to one aspect of the techniques described therein, an array of nucleic acid probes is fabricated at known locations on a substrate or chip. A fluorescently labeled nucleic acid is then brought into contact with the chip and a scanner generates an image file (which is processed into a cell file) indicating the locations where the labeled nucleic acids bound to the chip. Based upon the cell file and identities of the probes at specific locations, it becomes possible to extract information such as the monomer sequence of DNA or RNA. Such systems have been used to form, for example, arrays of DNA that may be used to study and detect mutations relevant to cystic fibrosis, the P53 gene (relevant to certain cancers), HIV, and other genetic characteristics.

[0004] Computer-aided techniques for monitoring gene expression using such arrays of probes have also been developed as disclosed in U.S. Patent Application No. 08/828,952 (Attorney Docket No. 16528X-028900US) and PCT Publication No. WO 97/10365 (Attorney Docket No. 16528X-017110PC), the contents of which are herein incorporated by reference. Many disease states are characterized by differences in the expression levels of various genes either through changes in the copy number of the genetic DNA or through changes in levels of transcription (*e.g.*, through control of initiation, provision of RNA precursors, RNA processing, etc.) of particular genes. For example, losses and gains of genetic material play an important role in malignant transformation and progression. Furthermore, changes in the expression (transcription) levels of particular genes (*e.g.*, oncogenes or tumor suppressors), serve as signposts for the presence and progression of various cancers.

[0005] It is desirable to identify genes having expression levels relevant to diagnosis of a diseased state by analyzing the expression levels of large numbers of genes in both diseased and normal individuals. Methods for collecting the expression level information have been developed. However, the user interfaces for gene ex-

pression monitoring systems that have been developed until now are designed to clearly present the expression of particular pre-selected genes. A user seeking to identify, *e.g.*, an oncogene or a tumor suppressor gene, must individually review the expression level of large numbers of genes and compare the expression levels between diseased and normal individuals. What is needed is a user interface that takes advantage of collected gene expression information to help the user to identify particular genes of interest.

SUMMARY OF THE INVENTION

[0006] The present invention provides innovative systems and methods for visualizing information collected from analyzing samples. The samples may include nucleic acids, proteins, or other polymers. Gene expression level as determined from analysis of a nucleic acid sample is one possible analysis result that may be visualized. In one embodiment, a computer system may display the expression levels of multiple genes simultaneously in a way that facilitates user identification of genes whose expression is significant to a characteristic such as disease or resistance to disease. Additionally, the computer system may facilitate display of further information about relevant genes once they are identified.

[0007] A first aspect of the invention provides a computer-implemented method for presenting expression level information as collected from first and second samples. The method includes steps of: displaying a first axis corresponding to expression level in the first sample, and displaying a second axis substantially perpendicular to the first axis, the second axis corresponding to expression level in the second sample. The method further includes a step of: for a selected expressed sequence, displaying a mark at a position. The position is selected relative to the first axis in accordance with an expression level of the selected expressed sequence in the first sample and relative to the second axis in accordance with an expression level of the selected expressed sequence in the second sample. A particularly useful application is displaying many marks simultaneously for many selected genes to discover which ones of the selected genes may be relevant to the characteristic.

[0008] A second aspect of the invention provides a computer-implemented method of presenting sample analysis information. The method includes steps of: displaying a first axis corresponding to a concentration of a compound in a first sample as determined by monitoring binding of the compound to a selected polymer having binding affinity to the compound, and displaying a second axis substantially perpendicular to the first axis. The second axis corresponds to a concentration of the compound in the second sample as determined by monitoring binding of the compound to the selected polymer. The method further preferably includes a step of displaying a mark at a position. The position is selected relative to the first axis in accordance with the concentration in

the first sample and relative to the second axis in accordance with the concentration in the second sample.

[0009] A further understanding of the nature and advantages of the inventions herein may be realized by reference to the remaining portions of the specification and the attached drawings.

BRIEF DESCRIPTION OF THE DRAWINGS

[0010] Fig. 1 illustrates an example of a computer system that may be used to execute software embodiments of the present invention.

[0011] Fig. 2 shows a system block diagram of a typical computer system.

[0012] Fig. 3 illustrates an overall system for forming and analyzing arrays of polymers including biological materials such as DNA or RNA.

[0013] Fig. 4 is an illustration of an embodiment of software for the overall system.

[0014] Fig. 5 shows a flowchart of a process of monitoring the expression of a gene by comparing hybridization intensities of pairs of perfect match and mismatch probes.

[0015] Fig. 6 shows a screen display illustrating gene expression levels for multiple genes as collected from both normal and diseased tissue.

[0016] Figs. 7A-7B show screen displays illustrating information about a particular gene selected from the display of Fig. 6.

DESCRIPTION OF SPECIFIC EMBODIMENTS

[0017] The present invention provides innovative methods of monitoring visualizing gene expression. In the description that follows, the invention will be described in reference to preferred embodiments. However, the description is provided for purposes of illustration and not for limiting the spirit and scope of the invention.

[0018] Fig. 1 illustrates an example of a computer system that may be used to execute software embodiments of the present invention. Fig. 1 shows a computer system 1 which includes a monitor 3, screen 5, cabinet 7, keyboard 9, and mouse 11. Mouse 11 may have one or more buttons such as mouse buttons 13. Cabinet 7 houses a CD-ROM drive 15 and a hard drive (not shown) that may be utilized to store and retrieve software programs including computer code incorporating the present invention. Although a CD-ROM 17 is shown as the computer readable medium, other computer readable media including floppy disks, DRAM, hard drives, flash memory, tape, and the like may be utilized. Cabinet 7 also houses familiar computer components (not shown) such as a processor, memory, and the like.

[0019] Fig. 2 shows a system block diagram of computer system 1 used to execute software embodiments of the present invention. As in Fig. 1, computer system 1 includes monitor 3 and keyboard 9. Computer system 1 further includes subsystems such as a central proc-

essor 50, system memory 52, I/O controller 54, display adapter 56, removable disk 58, fixed disk 60, network interface 62, and speaker 64. Removable disk 58 is representative of removable computer readable media like 5 floppies, tape, CD-ROM, removable hard drive, flash memory, and the like. Fixed disk 60 is representative of an internal hard drive or the like. Other computer systems suitable for use with the present invention may include additional or fewer subsystems. For example, another computer system could include more than one processor 50 (*i.e.*, a multi-processor system) or memory cache.

[0020] Arrows such as 66 represent the system bus architecture of computer system 1. However, these arrows are illustrative of any interconnection scheme serving to link the subsystems. For example, display adapter 56 may be connected to central processor 50 through a local bus or the system may include a memory cache. Computer system 1 shown in Fig. 2 is but an example of a computer system suitable for use with the present invention. Other configurations of subsystems suitable for use with the present invention will be readily apparent to one of ordinary skill in the art. In one embodiment, the computer system is an IBM compatible personal computer.

[0021] The VILSIPS™ and GeneChip™ technologies provide methods of making and using very large arrays of polymers, such as nucleic acids, on very small chips. See U.S. Patent No. 5,143,854 and PCT Patent Publication Nos. WO 90/15070 and 92/10092, each of which is hereby incorporated by reference for all purposes. Nucleic acid probes on the chip are used to detect complementary nucleic acid sequences in a sample nucleic acid of interest (the "target" nucleic acid).

[0022] It should be understood that the probes need not be nucleic acid probes but may also be other receptors, such as antibodies, or polymers such as peptides. Peptide probes may be used to detect the concentration of other peptides, proteins, or other compounds in a sample. The probes must be carefully selected to have bonding affinity to the compound whose concentration they are to be used to measure.

[0023] In one embodiment, the present invention provides methods of visualizing information relating to the concentration of compounds in a sample as measured by monitoring affinity of the compounds to probes. In a particular application, the concentration information is generated by analysis of hybridization intensity files for a chip containing hybridized nucleic acid probes. The hybridization of a nucleic acid sample to certain probes may represent the expression level of one more genes or expressed sequence tags (ESTs). The expression level of a gene or EST is herein understood to be the concentration within a sample of mRNA or protein that would result from the transcription of the gene or EST.

[0024] Expression level information visualized by virtue of the present invention need not be obtained from probes but may originate from any source. If the expres-

sion information is collected from a probe array, the probe array need not meet any particular criteria for size and density. Furthermore, the present invention is not limited to visualizing fluorescent measurements of bondings such as hybridizations but may be readily utilized to visualize other measurements.

[0025] Concentration of compounds other than nucleic acids may be visualized according to one embodiment of the present invention. For example, a probe array may include peptide probes which may be exposed to protein samples, polypeptide samples, or other compounds which may or may not bond to the peptide probes. By appropriate selection of the peptide probes, one may detect the presence or absence of particular compounds which would bond to the peptide probes.

[0026] For purposes of illustration, the present invention is described as being part of a system that designs a chip mask, synthesizes the probes on the chip, labels nucleic acids from a target sample, and scans the hybridized probes. Such a system is set forth in U.S. Patent No. 5,571,639 which is hereby incorporated by reference for all purposes. However, the present invention may be used separately from the overall system for analyzing data generated by such systems, such as at remote locations, or for visualizing the results of other systems for generating expression information, or for visualizing concentrations of polymers other than nucleic acids.

[0027] Fig. 3 illustrates a computerized system for forming and analyzing arrays of biological materials such as RNA or DNA. A computer 100 is used to design arrays of biological polymers such as RNA or DNA. The computer 100 may be, for example, an appropriately programmed IBM personal computer compatible running Windows NT including appropriate memory and a CPU as shown in Figs. 1 and 2. The computer system 100 obtains inputs from a user regarding characteristics of a gene of interest, and other inputs regarding the desired features of the array. Optionally, the computer system may obtain information regarding a specific genetic sequence of interest from an external or internal database 102 such as GenBank. The output of the computer system 100 is a set of chip design computer files 104 in the form of, for example, a switch matrix, as described in PCT application WO 92/10092, and other associated computer files.

[0028] The chip design files are provided to a system 106 that designs the lithographic masks used in the fabrication of arrays of molecules such as DNA. The system or process 106 may include the hardware necessary to manufacture masks 110 and also the necessary computer hardware and software 108 necessary to lay the mask patterns out on the mask in an efficient manner. As with the other features in Fig. 3, such equipment may or may not be located at the same physical site, but is shown together for ease of illustration in Fig. 3. The system 106 generates masks 110 or other synthesis patterns such as chrome-on-glass masks for use in the fab-

rication of polymer arrays.

[0029] The masks 110, as well as selected information relating to the design of the chips from system 100, are used in a synthesis system 112. Synthesis system 112 includes the necessary hardware and software used to fabricate arrays of polymers on a substrate or chip 114. For example, synthesizer 112 includes a light source 116 and a chemical flow cell 118 on which the substrate or chip 114 is placed. Mask 110 is placed between the light source and the substrate/chip, and the two are translated relative to each other at appropriate times for deprotection of selected regions of the chip. Selected chemical reagents are directed through flow cell 118 for coupling to deprotected regions, as well as for washing and other operations. All operations are preferably directed by an appropriately programmed computer 119, which may or may not be the same computer as the computer(s) used in mask design and mask making.

[0030] The substrates fabricated by synthesis system 112 are optionally diced into smaller chips and exposed to marked targets. The targets may or may not be complementary to one or more of the molecules on the substrate. The targets are marked with a label such as a fluorescein label (indicated by an asterisk in Fig. 3) and placed in scanning system 120. Scanning system 120 again operates under the direction of an appropriately programmed digital computer 122, which also may or may not be the same computer as the computers used in synthesis, mask making, and mask design. The scanner 120 includes a detection device 124 such as a confocal microscope or CCD (charge-coupled device) that is used to detect the location where labeled target has bound to the substrate. The output of scanner 120 is an image file(s) 124 indicating, in the case of fluorescein labeled target, the fluorescence intensity (photon counts or other related measurements, such as voltage) as a function of position on the substrate. Since higher photon counts will be observed where the labeled target has bound more strongly to the array of polymers, and since the monomer sequence of the polymers on the substrate is known as a function of position, it becomes possible to determine the sequence(s) of polymer(s) on the substrate that are complementary to the target.

[0031] The image file 124 is provided as input to an analysis system 126 that incorporates the visualization and analysis methods of the present invention. Again, the analysis system may be any one of a wide variety of computer system. The present invention provides various methods of analyzing and visualizing the chip design files and the image files, providing appropriate output 128. The chip design need not include any particular number of probes. It should be understood that the present invention does not require any particular source of expression level information.

[0032] Fig. 4 provides a simplified illustration of the overall software system used in the operation of one embodiment of the invention. As shown in Fig. 4, the sys-

tem first identifies the nucleotide sequence(s) or targets that would be of interest in a particular expression level analysis at step 202. The sequences of interest correspond to mRNA transcripts of one or more genes, ESTs or nucleic acids derived from the mRNA transcripts. Sequence selection may be provided via manual input of text files or may be from external sources such as GenBank.

[0033] At step 204 the system evaluates the sequences of interest to determine or assist the user in determining which probes would be desirable on the chip, and provides an appropriate "layout" on the chip for the probes. The process of selecting probes for an expression level analysis is explained in PCT Publication No. WO 97/10365, the contents of which are herein incorporated by reference. An alternative probe selection process that does not require prior knowledge of sequences of interest is explained in PCT Publication No. WO97/27317 (Attorney Docket No. 18547-019410PC), the contents of which are herein incorporated by reference. Further general background on probe selection is found in PCT Publication No. WO95/11995 (Attorney Docket No. 18547-004111PC) and PCT Publication No. WO97/29212 (Attorney Docket No. 18547-018540PC), the contents of which are herein incorporated by reference. The term "perfect match probe" refers to a probe that has a sequence that is perfectly complementary to a particular target sequence. The test probe is typically perfectly complementary to a portion (subsequence) of the target sequence. The term "mismatch control" or "mismatch probe" refer to probes whose sequence is deliberately selected not to be perfectly complementary to a particular target sequence. For each mismatch (MM) control in an array there typically exists a corresponding perfect match (PM) probe that is perfectly complementary to the same particular target sequence.

[0034] The process compares hybridization intensities of pairs of perfect match and mismatch probes that are preferably covalently attached to the surface of a substrate or chip. Most preferably, the nucleic acid probes have a density greater than about 60 different nucleic acid probes per 1 cm² of the substrate.

[0035] Initially, nucleic acid probes are selected that are complementary to the target sequence. These probes are the perfect match probes. Another set of probes is specified that are intended to be not perfectly complementary to the target sequence. These probes are the mismatch probes and each mismatch probe includes at least one nucleotide mismatch from a perfect match probe. Accordingly, a mismatch probe and the perfect match probe to which it is identical except for one base make up a pair. As mentioned earlier, the nucleotide mismatch is preferably near the center of the mismatch probe.

[0036] The probe lengths of the perfect match probes are typically chosen to exhibit detectably greater hybridization with the target sequence relative to the mismatch probes. For example, the nucleic acid probes may be

all 20-mers. However, probes of varying lengths may also be synthesized on the substrate for any number of reasons including resolving ambiguities.

[0037] Again referring to Fig. 4, at step 206 the masks for the synthesis are designed. At step 208 the software utilizes the mask design and layout information to make the DNA or other polymer chips. This step 208 will control, among other things, relative translation of a substrate and the mask, the flow of desired reagents through a flow cell, the synthesis temperature of the flow cell, and other parameters. At step 210, another piece of software is used in scanning a chip thus synthesized and exposed to a labeled target. The software controls the scanning of the chip, and stores the data thus obtained in a file that may later be utilized to extract hybridization information.

[0038] At step 212 a computer system utilizes the layout information and the fluorescence information to evaluate the hybridized nucleic acid probes on the chip. Among the important pieces of information obtained from DNA chips are the relative fluorescent intensities obtained from the perfect match probes and mismatch probes. These intensity levels are used to estimate an expression level for a gene or EST. The computer system used for analysis will preferably have available other details of the experiment including possibly the gene name, gene sequence, probe sequences, probe locations on the substrate, and the like.

[0039] According to the present invention, at step 214, the same computer system used for analysis or another one displays the expression level information in a format useful for identifying genes of interest. The visualized expression level information may include information collected from multiple applications of one or more previous steps of Fig. 4.

[0040] Fig. 5 is a flowchart describing steps of estimating an expression level for a particular gene and determining whether the expression level is sufficiently high to be displayed. At step 952, the computer system receives raw scan data of N pairs of perfect match and mismatch probes. In a preferred embodiment, the hybridization intensities are photon counts from a fluorescently labeled target that has hybridized to the probes on the substrate. For simplicity, the hybridization intensity of a perfect match probe will be designed "I_{pm}" and the hybridization intensity of a mismatch probe will be designed "I_{mm}".

[0041] Hybridization intensities for a pair of probes are retrieved at step 954. The background signal intensity is subtracted from each of the hybridization intensities of the pair at step 956. Background subtraction can also be performed on all the raw scan data at the same time.

[0042] At step 958, the hybridization intensities of the pair of probes are compared to a difference threshold (D) and a ratio threshold (R). It is determined if the difference between the hybridization intensities of the pair (I_{pm} - I_{mm}) is greater than or equal to the difference threshold AND the quotient of the hybridization intensities

ties of the pair (I_{pm} / I_{mm}) is greater than or equal to the ratio threshold. The difference thresholds are typically user defined values that have been determined to produce accurate expression monitoring of a gene or genes. In one embodiment, the difference threshold is 20 and the ratio threshold is 1.2.

[0043] If $I_{pm} - I_{mm} > D$ and $I_{pm} / I_{mm} > R$, the value NPOS is incremented at step 960. In general, NPOS is a value that indicates the number of pairs of probes which have hybridization intensities indicating that the gene is likely expressed. NPOS is utilized in a determination of the expression of the gene.

[0044] At step 962, it is determined if $I_{mm} - I_{pm} > D$ and $I_{mm} / I_{pm} > R$. If these expressions are true, the value NNEG is incremented at step 964. In general, NNEG is a value that indicates the number of pairs of probes which have hybridization intensities indicating that the gene is likely not expressed. NNEG, like NPOS, is utilized in a determination of the expression of the gene.

[0045] For each pair that exhibits hybridization intensities either indicating the gene is expressed or not expressed, a log ratio value (LR) and intensity difference value (IDIF) are calculated at step 966. LR is calculated by the log of the quotient of the hybridization intensities of the pair (I_{pm} / I_{mm}). The IDIF is calculated by the difference between the hybridization intensities of the pair ($I_{pm} - I_{mm}$). If there is a next pair of hybridization intensities at step 968, they are retrieved at step 954.

[0046] At step 972, a decision matrix is utilized to indicate if the gene is expressed. The decision matrix utilizes the values N, NPOS, NNEG, LR (multiple LRs), and IDIF (multiple IDIFs). The following four assignments are performed:

$$P1 = NPOS / NNEG$$

$$P2 = NPOS / N$$

$$P3 = \text{SUM}(\text{LR}) / N$$

$$P4 = \text{SUM}(\text{IDIF}) / N$$

These P values are then utilized to determine if the gene is expressed and if the expression level should be displayed. In a preferred embodiment, the expression level of a gene should be displayed if:

$$P1 > 2.2$$

$$P2 > 0.3$$

$$P3 > 0.8$$

$$P4 > 30$$

[0047] Once all the pairs of probes have been processed and the expression of the gene indicated, an average of the IDIF values for the probes that incremented

NPOS or NNEG is calculated at step 975, which is utilized as an expression level. Of course, other values including one of P1 through P4 could be used to indicate expression level.

[0048] For simplicity, Fig. 5 was described in reference to a single gene or EST. However, the visualization system of the present invention displays expression results for many genes to facilitate discovery of genes of interest or ESTs. Furthermore, the present invention contemplates display of expression levels of a single gene or ESTs' as collected from two or more different samples such as tissue samples. The sample sources preferably differ in some characteristic. It will be understood that when the term "sample" is used herein, measurements made on a single "sample" can be based on an aggregation of multiple sample collection events or even multiple organisms.

[0049] Fig. 6 shows a screen display illustrating gene expression levels for multiple genes as collected from two tissue samples. A displayed horizontal axis 1002 represents expression level measured in one or more nucleic acid samples taken from the first tissue sample. A displayed vertical axis 1004 represents expression level in one or more nucleic acid samples taken from the second tissue sample. Each of marks 1006 represent a particular gene whose expression level has been measured in both the first and second tissue samples. Each mark 1006 is placed at a distance from vertical axis 1004 corresponding to expression level in the first tissue sample and at a distance from the horizontal axis 1002 corresponding to expression level in the second tissue sample.

[0050] The expression levels used for determining the position of marks 1006 are preferably taken from the result of step 975. The position of each of marks 1006 depends on two iterations of the steps of Fig. 5, once for the sample taken from the first tissue sample and once for the sample taken from the second tissue sample. However, a mark is preferably displayed only if one of the samples meets the threshold criteria at step 972.

[0051] In the depicted representative screen display, the first tissue sample is a cancerous tissue sample and the second tissue sample is a normal tissue sample. The individual marks represent the expression levels of selected genes in both cancerous and normal tissue. A first group of marks 1008 represent genes that are neither tumor suppressors nor oncogenes since their expression levels are roughly similar for both normal and cancerous tissue. These marks 1008 fall roughly along a line which is rotated 45 degrees from each of the axes. A second group of marks 1010 represent genes that are likely oncogenes since their expression levels are found to be significantly higher in cancerous tissue than in normal tissue. A third group of marks 1012 represent genes that are likely tumor suppressors since their expression levels are found to be significantly higher in normal tissue than in cancerous tissue. It will be appreciated that expression levels for large numbers of genes can be re-

viewed at once to discover the oncogenes and tumor suppressors.

[0052] Although in the depicted display, the two types of tissue are normal tissue and cancerous tissue, the present invention would aid in the discovery of genes whose expression is associated with any characteristic that varies among tissue samples. For example, one can compare expression results from tissue from individuals who have been exposed to HIV but remain infected to tissue obtained from infected individuals to identify genes conferring resistance to HIV. One can compare expression results between tissue from plants that survive drought to plants that do not. One can compare expression levels among tissue samples at successive stages or severity levels of the same disease, among tissue samples where different ultimate outcomes of the disease (e.g., patient death or remission) are known, among diseased tissue samples that have been subject to different treatment regimes including e.g., chemotherapy, antisense RNA, etc. For cancers, one can compare expression levels between malignant cells and non-malignant cells. Also expression levels can be compared among different organs, between species, and among different stages of development of an organ.

[0053] It will be appreciated that the present invention also encompasses displays with more than two dimensions. A third visual dimension can be used to illustrate expression level from a third tissue sample. The time dimension can also be used to illustrate successive groups of two or three tissue samples at successive time periods. The time dimension can be also used to correspond to tissue samples obtained at, e.g., successive stages of a disease.

[0054] Other interface methods corresponding to human senses other than sight can also be incorporated within the presentation system of the present invention. The senses may correspond to additional dimensions. For example, marks can be displayed in succession accompanied by a sound having characteristics corresponding to expression level in another tissue sample.

[0055] The user can employ a cursor 1014 to identify a particular mark as being of interest. Cursor 1014 can be moved to a particular mark by use of, e.g., mouse 11. Once cursor 1014 is over a mark of interest, the mark can be selected by, e.g., depression of one of mouse buttons 13. Selection of a particular mark can be facilitated by use of a zoom display feature (not shown). Once a particular mark is selected, further information is displayed about the gene represented by the mark. A special mouse can transmit a tactile sensation back to the user corresponding to expression level in a tissue sample as the user passes the mouse over a corresponding mark.

[0056] It will be appreciated that the display of Fig. 6 is not limited to expression information. The two dimensions of Fig. 6 may correspond to indicators of the presence of various polymers other than nucleic acids in two different samples. For example, each mark may corre-

spond to a different polymer, polypeptide, or other compound. The distance of the mark from each axis would correspond to a measure of presence of the particular polymer in the sample corresponding to the axis. One possible measure is produced by fluorescently tagging polymer samples such as protein samples and exposing a probe array such as a peptide probe array to the protein samples. The fluorescent intensity of the probes will then correspond to the bonding affinity of the sample to the probes. The intensity measurement or a measurement derived from the intensity measurement may then be used to position the marks of Fig. 6.

[0057] Fig. 7A shows a screen display giving information about a particular gene selected from the display of Fig. 6. A cluster number 702, a GenBank accession number 704, and a verbal description 706 for the selected gene are displayed. The user can also select a number of marks 1006 by circling them with cursor 1014. Then a list of information as shown in Fig. 7A is displayed for all the genes corresponding to the selected marks.

[0058] By selecting GenBank accession number 704 with another cursor (not shown), the user can direct retrieval of the GenBank information for the selected gene. If the GenBank information is not available locally, the retrieval process can include formulating a query and transmitting the query to a GenBank web site. Once the GenBank information is retrieved, it can also be displayed. Fig. 7B depicts the GenBank information for the gene identified in Fig. 7A.

[0059] In the foregoing specification, the invention has been described with reference to specific exemplary embodiments thereof. It will, however, be evident that various modifications and changes may be made thereunto without departing from the broader spirit and scope of the invention as set forth in the appended claims and their full scope of equivalents.

Claims

1. A computer-implemented method of presenting expression level information as collected from first and second samples, said method comprising the steps of:

displaying a first axis corresponding to expression level in said first sample;

displaying a second axis substantially perpendicular to said first axis, said second axis corresponding to expression level in said second sample; and

for a selected expressed sequence, displaying a mark at a position, wherein said position is selected relative to said first axis in accordance with an expression level of said selected expressed sequence in said first sample and relative to said second axis in accordance with an

- expression level of said selected expressed sequence in said second sample.
2. The method of claim 1 wherein said selected expressed sequence comprises a gene. 5
 3. The method of claim 1 wherein said selected expressed sequence comprises a portion of a gene.
 4. The method of claim 1 further comprising the step of repeating said displaying a mark step for a plurality of selected expressed sequences. 10
 5. The method of claim 1 further comprising the steps of: 15
 - monitoring said expression level of said expressed sequence in said first sample and said second sample.
 6. The method of claim 3 wherein said monitoring step for one of said samples comprises substeps of: 20
 - inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, said perfect match probes being perfectly complementary to a target nucleic acid sequence indicative of expression of said selected gene and said mismatch probes having at least one base mismatch with said target sequence, and said hybridization intensities indicating hybridization affinity between said perfect match and mismatch probes and a sample nucleic acid sequence from said one of said samples; 25
 - comparing the hybridization intensities of each pair of perfect match probe and mismatch probe; and 30
 - generating said expression level for said expressed sequence and said one of said samples responsive to results of said comparing step. 35
 7. The method of claim 6 further comprising the step of: 40
 - comparing a difference between hybridization intensities of perfect match and mismatch probes at a base position to a difference threshold. 45
 8. The method of claim 7 further comprising the step of: 50
 - comparing a quotient of hybridization intensities of perfect match and mismatch probes at a base position to a ratio threshold.
 9. The method of claim 6 further comprising the steps of: 55
 - a) counting a probe pair as a positive probe pair to increment a positive probe pair count if a perfect match probe intensity minus a mismatch probe intensity exceeds a difference threshold and said perfect match probe intensity divided by said mismatch probe intensity exceeds a ratio threshold; 60
 - b) counting said probe pair as a negative probe pair to increment a negative probe pair count if said mismatch probe intensity minus said perfect match probe intensity exceeds said difference threshold and said mismatch probe intensity divided by said perfect match probe intensity exceeds said ratio threshold; and 65
 - c) computing a logarithmic ratio of said perfect match probe intensity to said mismatch probe intensity.
 10. The method of claim 9 further comprising the steps of: 70
 - repeating said a), b), and c) steps for each of said probe pairs, accumulating a sum of differences of said perfect match and mismatch probe intensities for probe pairs that cause; and determining an expression level of said selected expressed sequence to be an average of said differences.
 11. The method of claim 1 further comprising the steps of: 75
 - receiving user input selecting said mark; and in response to said user input, displaying information about said selected expressed sequence.
 12. The method of claim 11 further comprising the steps of: 80
 - in response to said user input, displaying information about said selected expressed sequence.
 13. The method of claim 12 wherein said information about said selected expressed sequence comprises a GenBank accession number.
 14. The method of claim 12 wherein said information about said selected expressed sequence comprises a GenBank database record for said selected expressed sequence.
 15. The method of claim 1 wherein said first sample and said second sample are collected from tissue samples differing in a particular characteristic.
 16. The method of claim 15 wherein said particular characteristic comprises presence of disease.
 17. The method of claim 15 wherein said particular

characteristic comprises a treatment strategy for a disease.

18. The method of claim 1 wherein said particular characteristic is a stage of a disease.

19. The method of claim 1 further comprising the step of:

displaying a third axis substantially perpendicular to said first axis and to said second axis in a three-dimensional display environment wherein said position of said mark is further selected relative to said third axis in accordance with an expression level of said selected expressed sequence in a third sample.

20. A computer-implemented method of presenting sample analysis information comprising the steps of:

displaying a first axis corresponding to a concentration of a compound in a first sample as determined by monitoring binding of said compound to a selected polymer having binding affinity to said compound;

displaying a second axis substantially perpendicular to said first axis, said second axis corresponding to a concentration of said compound in said second sample as determined by monitoring binding of said compound to said selected polymer; and

displaying a mark at a position, wherein said position is selected relative to said first axis in accordance with said concentration in said first sample and relative to said second axis in accordance with said concentration in said second sample.

21. The method of claim 20 wherein said selected polymer comprises a nucleic acid sequence.

22. The method of claim 20 wherein said selected polymer comprises a protein.

23. The method of claim 21 further comprising the step of:

obtaining said concentration of said compound in said first sample by exposing said first sample to a plurality of nucleic acid probes.

24. The method of claim 22 further comprising the step of:

obtaining said concentration of said compound in said first sample by exposing said first sample to a plurality of peptide probes.

25. A computer program product for presenting expression level information as collected from first and

second samples, said product comprising::

code for displaying a first axis corresponding to expression level in said first sample;

code for displaying a second axis substantially perpendicular to said first axis, said second axis corresponding to expression level in said second sample;

code for, for a selected expressed sequence, displaying a mark at a position, wherein said position is selected relative to said first axis in accordance with an expression level of said selected expressed sequence in said first sample and relative to said second axis in accordance with an expression level of said selected expressed sequence in said second sample; and a computer-readable storage medium for storing the codes.

26. The product of claim 25 wherein said selected expressed sequence comprises a gene.

27. The product of claim 25 wherein said selected expressed sequence comprises a portion of a gene.

28. The product of claim 25 further comprising code for repeatedly applying said displaying a mark code for a plurality of selected expressed sequences.

29. The product of claim 25 further comprising: code for monitoring said expression level of said expressed sequence in said first sample and said second sample.

30. The product of claim 27 wherein said monitoring step for one of said samples comprises:

code for inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, said perfect match probes being perfectly complementary to a target nucleic acid sequence indicative of expression of said selected gene and said mismatch probes having at least one base mismatch with said target sequence, and said hybridization intensities indicating hybridization affinity between said perfect match and mismatch probes and a sample nucleic acid sequence from said one of said samples;

comparing the hybridization intensities of each pair of perfect match probe and mismatch probe; and

generating said expression level for said expressed sequence and said one of said samples responsive to results of said comparing step.

31. The product of claim 30 further comprising:

code for comparing a difference between hybridization intensities of perfect match and mismatch probes at a base position to a difference threshold.

32. The product of claim 31 further comprising:

code for comparing a quotient of hybridization intensities of perfect match and mismatch probes at a base position to a ratio threshold.

33. The product of claim 30 further comprising:

a) code for counting a probe pair as a positive probe pair to increment a positive probe pair count if a perfect match probe intensity minus a mismatch probe intensity exceeds a difference threshold and said perfect match probe intensity divided by said mismatch probe intensity exceeds a ratio threshold;

b) code for counting said probe pair as a negative probe pair to increment a negative probe pair count if said mismatch probe intensity minus said perfect match probe intensity exceeds said difference threshold and said mismatch probe intensity divided by said perfect match probe intensity exceeds said ratio threshold; and

c) code for computing a logarithmic ratio of said perfect match probe intensity to said mismatch probe intensity.

34. The product of claim 33 further comprising:

code for repeatedly applying said a), b), and c) codes for each of said probe pairs, accumulating a sum of differences of said perfect match and mismatch probe intensities for probe pairs that cause; and

code for determining an expression level of said selected expressed sequence to be an average of said differences.

35. The product of claim 25 further comprising:

code for receiving user input selecting said mark; and
code for, in response to said user input, displaying information about said selected expressed sequence.

36. The product of claim 35 further comprising:

code for, in response to said user input, displaying information about said selected expressed sequence.

37. The product of claim 36 wherein said information about said selected expressed sequence comprises a GenBank accession number.

38. The product of claim 36 wherein said information about said selected expressed sequence comprises a GenBank database record for said selected expressed sequence.

39. The product of claim 25 wherein said first sample and said second sample are collected from tissue samples differing in a particular characteristic.

40. The product of claim 39 wherein said particular characteristic comprises presence of disease.

41. The product of claim 39 wherein said particular characteristic comprises a treatment strategy for a disease.

42. The product of claim 25 wherein said particular characteristic is a stage of a disease.

43. The product of claim 25 further comprising the step of:

displaying a third axis substantially perpendicular to said first axis and to said second axis in a three-dimensional display environment wherein said position of said mark is further selected relative to said third axis in accordance with an expression level of said selected expressed sequence in a third sample.

44. A computer program product for presenting sample analysis information comprising:

code for displaying a first axis corresponding to a concentration of a compound in a first sample as determined by monitoring binding of said compound to a selected polymer having bonding affinity to said compound;

code for displaying a second axis substantially perpendicular to said first axis, said second axis corresponding to concentration of said compound in a second sample as determined by monitoring binding of said compound to said selected polymer;

code for displaying a mark at a position, wherein said position is selected relative to said first axis in accordance with said concentration in said first sample and relative to said second axis in accordance with said concentration in said second sample; and

a computer-readable storage medium that stores the codes.

45. The product of claim 44 wherein said selected polymer comprises a nucleic acid sequence.

46. The product of claim 44 wherein said selected polymer comprises a protein.

47. A computer system comprising a display, a processor, and a memory that stores instructions for configuring said processor to:

display a first axis corresponding to expression level in said first sample;
display a second axis substantially perpendicular to said first axis, said second axis corresponding to expression level in said second sample; and
for a selected expressed sequence, display a mark at a position, wherein said position is selected relative to said first axis in accordance with an expression level of said selected expressed sequence in said first sample and relative to said second axis in accordance with an expression level of said selected expressed sequence in said second sample.

48. A computer system comprising a display, a processor, and a memory that stores instructions for configuring said processor to:

display a first axis corresponding to a concentration of a compound in a first sample as determined by monitoring binding of said compound to a selected polymer having binding affinity to said compound;
display a second axis substantially perpendicular to said first axis, said second axis corresponding to a concentration of said compound in said second sample as determined by monitoring binding of said compound to said selected polymer; and
display a mark at a position, wherein said position is selected relative to said first axis in accordance with said concentration in said first sample and relative to said second axis in accordance with said concentration in said second sample.

49. A method of monitoring gene expression in first and second samples, the method comprising presenting expression level information as collected from said first and second samples, in accordance with any of claims 1 to 19, and using the displayed mark to monitor said gene expression.

50. A method of analysing first and second samples, the method comprising presenting sample analysis information relating to said first and second samples, in accordance with any of claims 20 to 24, and using the displayed mark to analyse said samples.

51. A method of identifying a gene of interest, the method comprising presenting expression level information relating to a gene, as collected from said first and second samples, in accordance with any of

claims 1 to 19, and using the displayed mark to identify whether said gene is of interest.

52. A method of identifying a gene having different expression levels in first and second samples, the method comprising presenting expression level information relating to said gene in accordance with any of claims 1 to 19, said expression level information being collected from said first and second samples, and using the displayed mark to identify said gene.

53. A method of identifying a gene having an effect on a characteristic of a tissue sample, the method comprising presenting expression level information relating to said gene in accordance with any of claims 1 to 19, said expression level information being collected from said first and second samples differing in said characteristic, and using the displayed mark to identify said gene.

54. A method of identifying a compound having different concentrations in first and second samples, the method comprising presenting sample analysis information relating to said first and second samples, in accordance with any of claims 20 to 24, and using the displayed mark to analyse said samples.

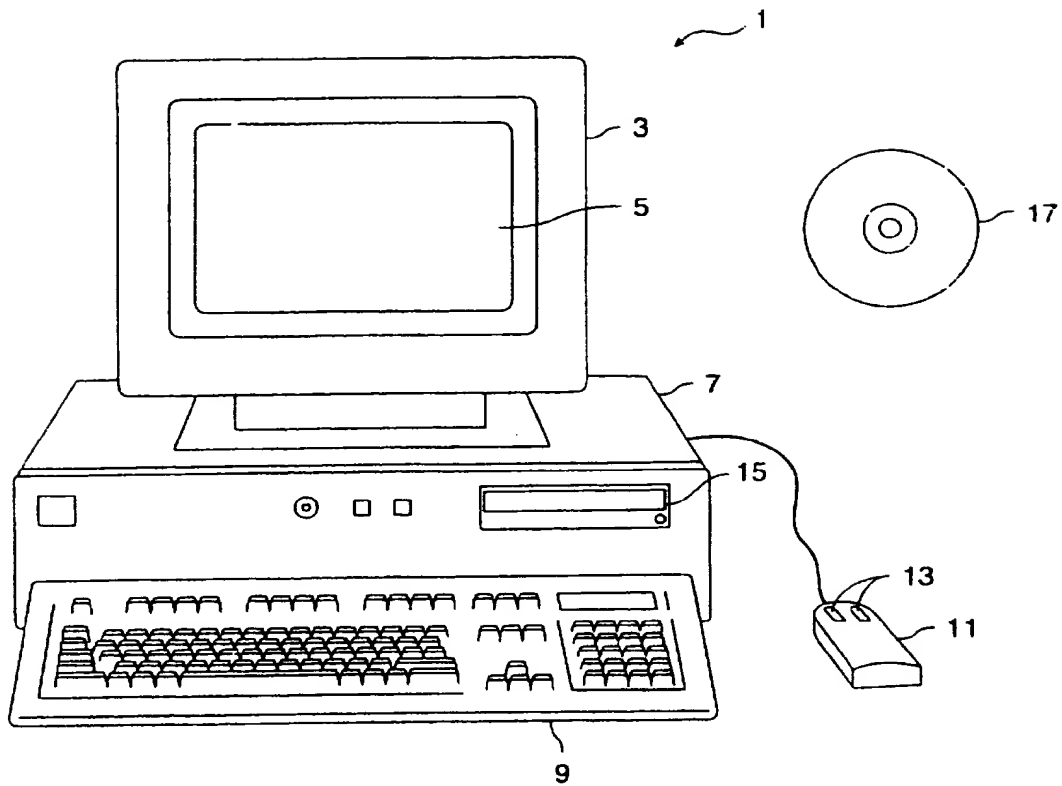


FIG. 1

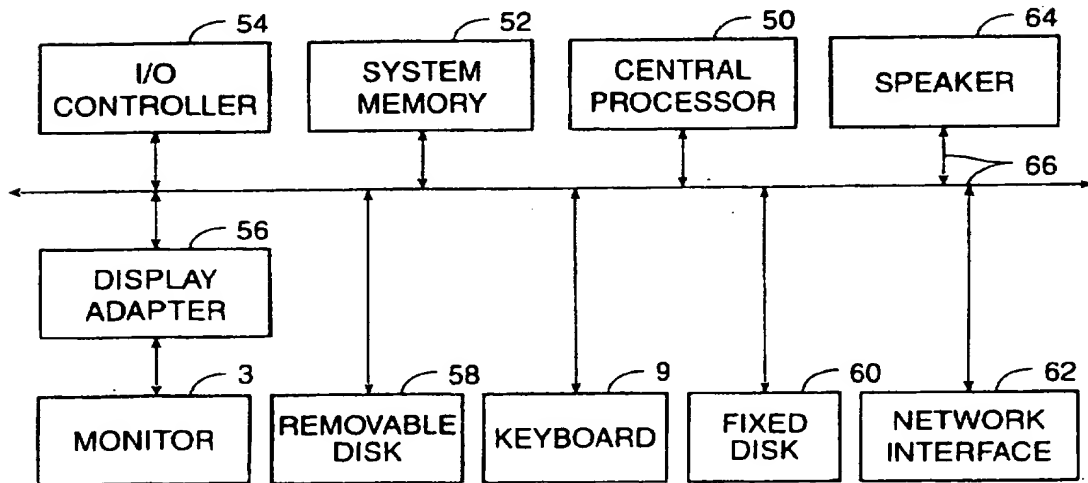


FIG. 2

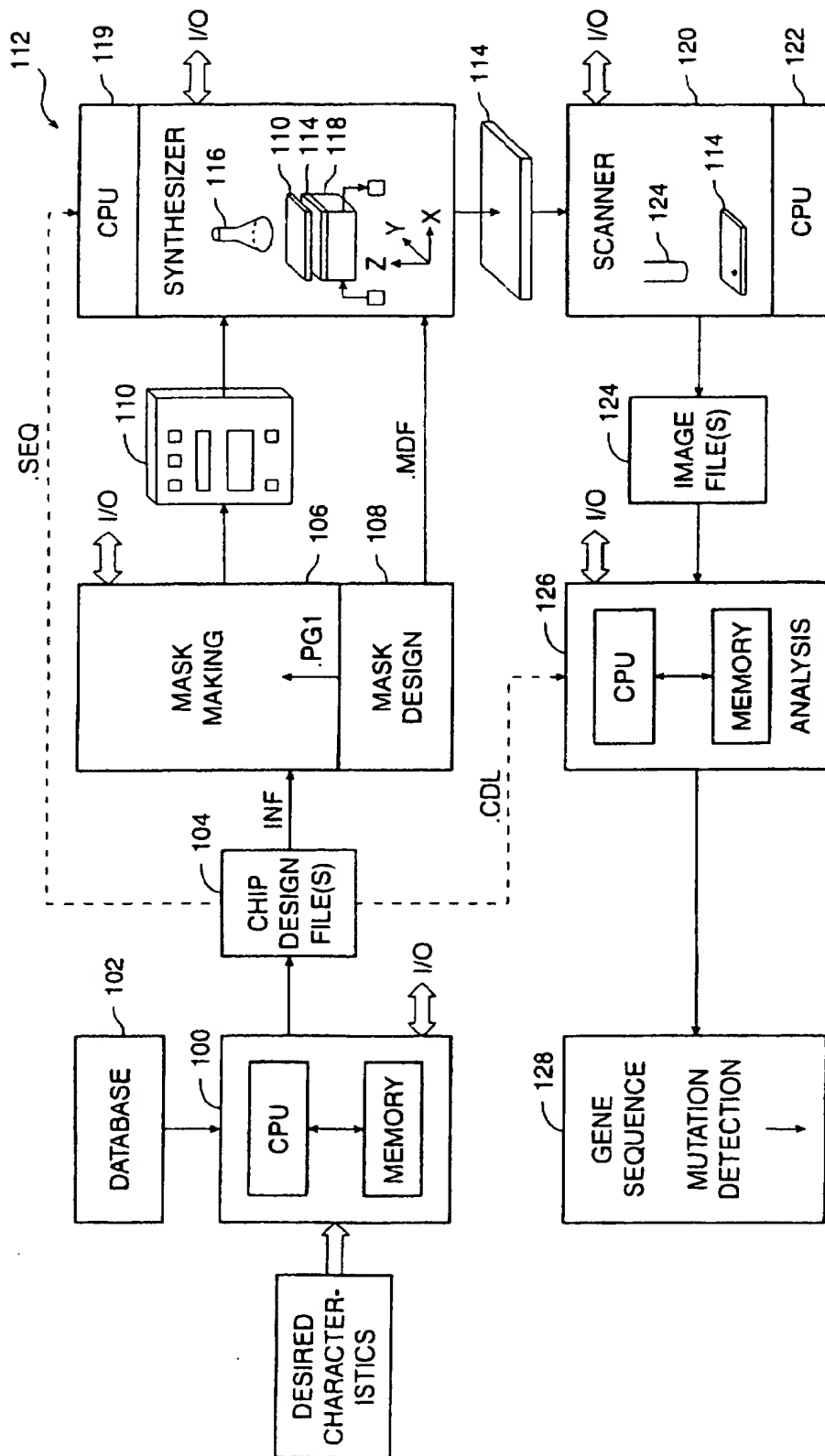


FIG. 3

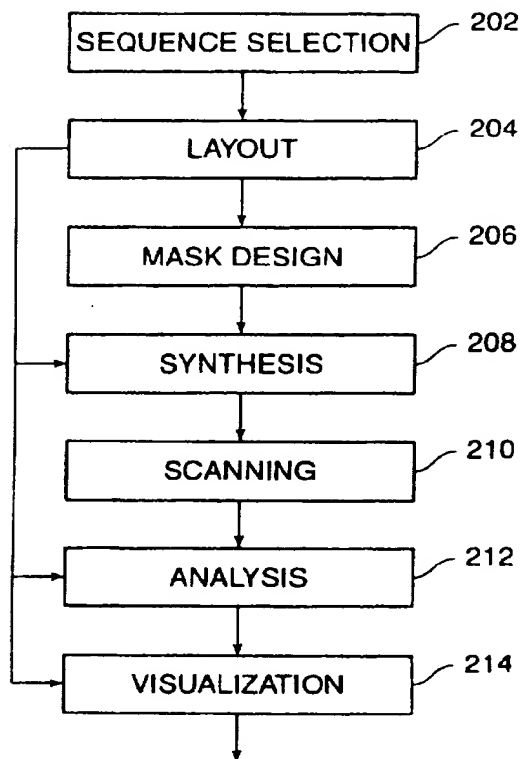


FIG. 4

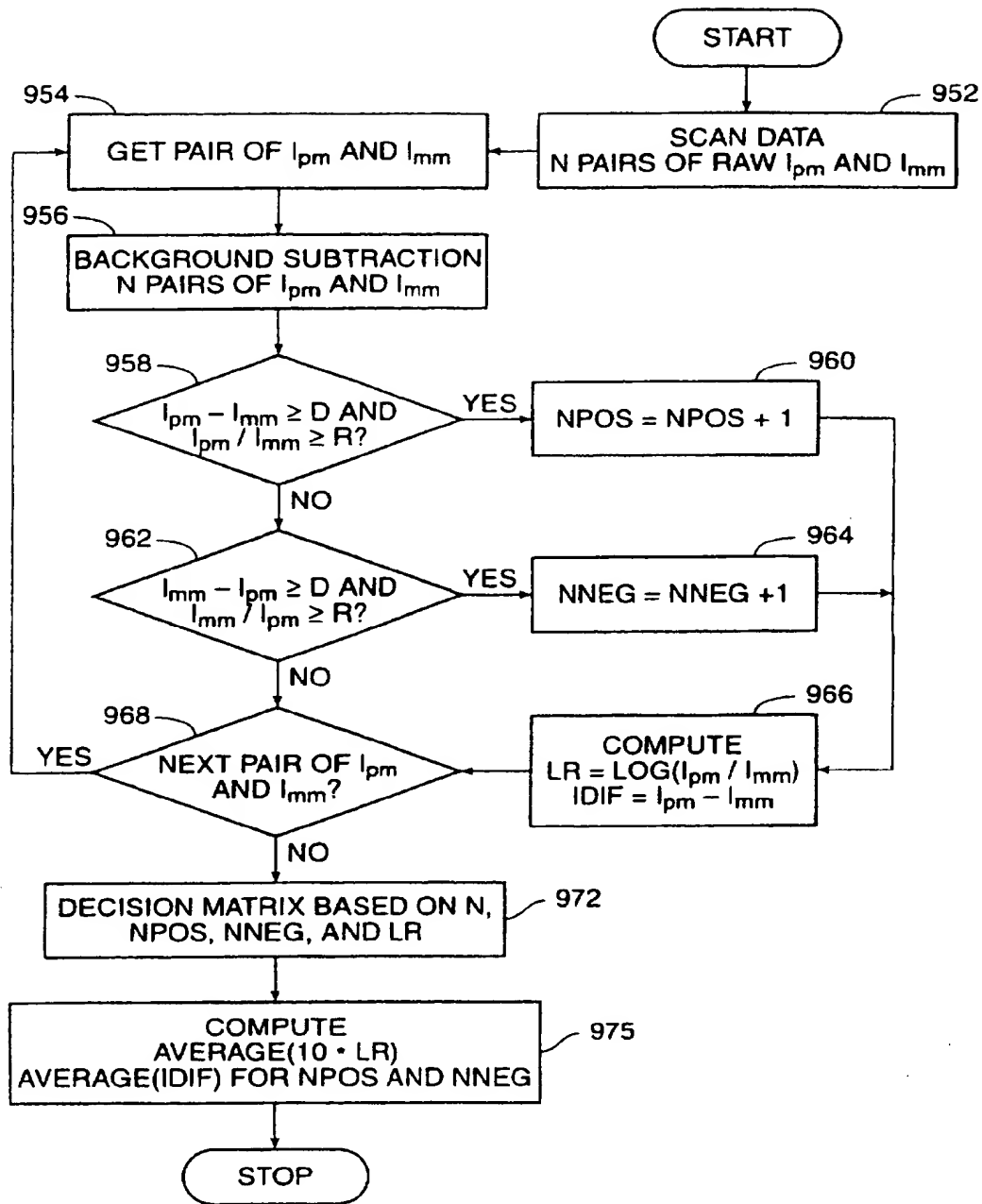


FIG. 5

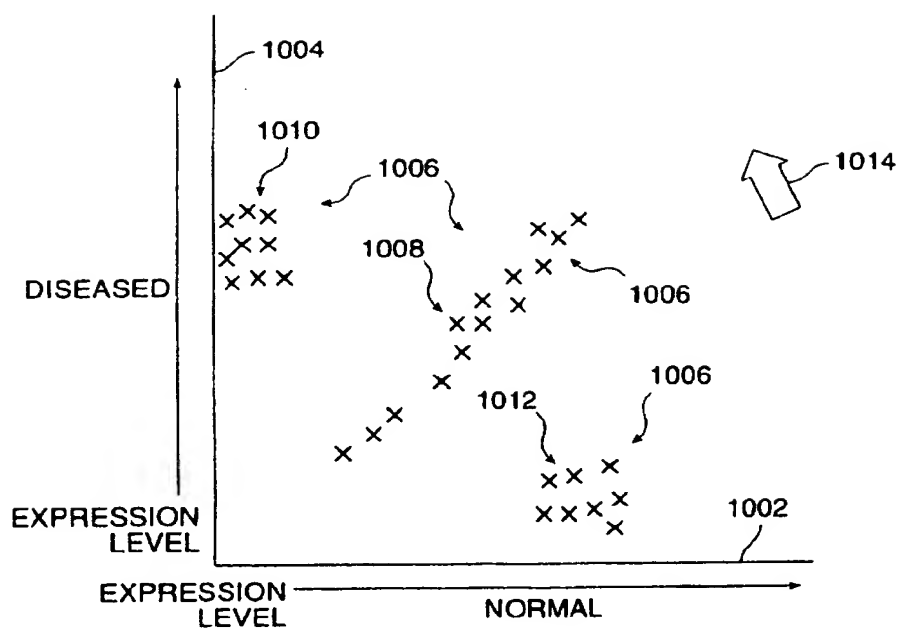


FIG. 6

CLUSTER	ACCESSION NUMBER	DESCRIPTION
Hsa.35	D11327	Human mRNA for protein-tyrosine phosphatase; complete cds.
702	704	706

FIG. 7A

http://www.ncbi.nlm.nih.gov/irx/cgi-bin/birx_doc?genbank+68169

LOCUS HUMLCPTP 2691 bp mRNA PRI 06-NOV-1992
 DEFINITION Human mRNA for protein-tyrosine phosphatase, complete cds.
 ACCESSION D11327
 NID g219901
 KEYWORDS protein-tyrosine phosphatase.
 SOURCE Human T cell, lambda-gt10 library, cDNA to mRNA.
 ORGANISM Homo sapiens
 Eukaryotae; mitochondrial eukaryotes; Metazoa; Chordata;
 Vertebrata; Mammalia; Eutheria; Primates; Catarrhini; Hominidae;
 Homo.
 REFERENCE 1 (bases 1 to 2691)
 AUTHORS Adachi,T., Sekiya,M., Isobe,M., Kumura,Y., Ogita,Z., Hinoda,Y.,
 Imai,K. and Yachi,A.
 TITLE Molecular cloning and chromosomal mapping of a human
 protein-tyrosine phosphatase LC-PTP
 JOURNAL Biochemical and Biophysical Research Communication 186, 1607-1615
 (1992)
 COMMENT Submitted (22-MAY-1992) to DDBJ by:
 Masaaki Adachi
 Sapporo Medical College
 S1W16 Chuo-ku
 Sapporo 060
 Japan
 Phone: 011-611-2111
 Fax: 011-613-1141.
 FEATURES
 source Location/Qualifiers
 1..2691
 /organism='Homo sapiens'
 /db_xref='taxon:9606'
 /cell_type='T cell'
 /clone_lib='lambda-gt10'
 gene 105..1187
 /gene='LC-PTP'
 CDS 105..1187
 /gene='LC-PTP'
 /codon_start=1
 /product='protein-tyrosine phosphatase'
 /db_xref='PID:d1002425'
 /db_xref='PID:g219902'
 /translation='MVQAHGGRSRAQPLTSLGAAMTQPPPEKTPAKKHVRLQERRGS
 NVALMLDVRSLGAVEPICSVNTPREVTLHFLRTAGHPLTRWALQRQPPSPKQLEEEFL
 KIPSNFVSPEDLDIPGHASKDRYKTIILPNPQSRVCLGRAQSQEDGDYINANYIRGYDG
 KEKVYIATQGPMPNTVSDFWEMVWQEEVSLIVMLTQLREGKEKCVHYWPTTEETYGPF
 QIRIQDMKECPEYTVRQLTIQYQEERRSVKHILFSAWPDHOTPESAGPLLRRLVAEEVE
 SPETAHPGFIIVHCSAGIGRTGCFIATRIGCQQLKARGEVDILGIVCQLRLDRGGM
 QTDEQYQFLHHTLALYAGQLPEEPSP'
 misc_feature 444..1172
 /gene='LC-PTP'
 /note='single catalytic domain'
 polyA_signal 2667..2672
 polyA_site 2691
 BASE COUNT 652 a 816 c 707 g 516 t
 ORIGIN
 1 ggagacagac agacagctgg caagagggcag cctggggggcc acagctgctt cagcagacct
 61 catggctgag tgagcctccc ctggggcccag caccacacct cagcatgggc caagcccatg
 121 gggggcgctc cagagcacag ccgttgacct tgtctttggg ggcagccatg acccagcctc
 181 cgcctgaaaa aacgccagcc aagaagcatg tgcgactgca ggagaggcgg ggctccaatg
 241 tggctctgat gctggacgtt cggtccctgg gggccgtaga acccatctgc tctgtgaaca
 301 caccgccgga ggccacccta cactttctgc gcactgctgg acacccccctt acccgctggg
 361 ccccttcagc ccagccaccc agccccaagc aactgggaaga agaattcttg aagatccctt
 421 caaaccttgt cagcccccga gacctggaca tccctggcca cgctccaaag gaccgatata
 481 agaccatctt gccaaatccc cagagccgtg tctgtctagg ccgggcacag agccaggagg...

FIG. 7B

THIS PAGE BLANK (USPTO)

(19)



Europäisches Patentamt

European Patent Office

Office européen des brevets



(11)

EP 0 935 210 A3

(12)

EUROPEAN PATENT APPLICATION

(88) Date of publication A3:
20.11.2002 Bulletin 2002/47

(51) Int Cl.7: G06F 19/00

(43) Date of publication A2:
11.08.1999 Bulletin 1999/32

(21) Application number: 99300900.0

(22) Date of filing: 08.02.1999

(84) Designated Contracting States:
AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU
MC NL PT SE
Designated Extension States:
AL LT LV MK RO SI

(72) Inventor: Mack, David H.
Menlo Park, C.A. 94025 (US)

(74) Representative:
O'Connell, David Christopher et al
Haseltine Lake & Co.,
Imperial House,
15-19 Kingsway
London WC2B 6UD (GB)

(30) Priority: 09.02.1998 US 20743

(71) Applicant: Affymetrix, Inc. (a California
Corporation)
Santa Clara, CA 95051 (US)

(54) Computer aided visualisation of expression comparison

(57) Innovative systems and methods for visualizing information collected from analyzing samples are provided. The samples may include nucleic acids, proteins, or other polymers. Gene expression level as determined from analysis of a nucleic acid sample is one possible analysis result that may be visualized. In one embodi-

ment, a computer system may display the expression levels of multiple genes simultaneously in a way that facilitates user identification of genes whose expression is significant to a characteristic such as disease or resistance to disease. Additionally, the computer system may facilitate display of further information about relevant genes once they are identified.

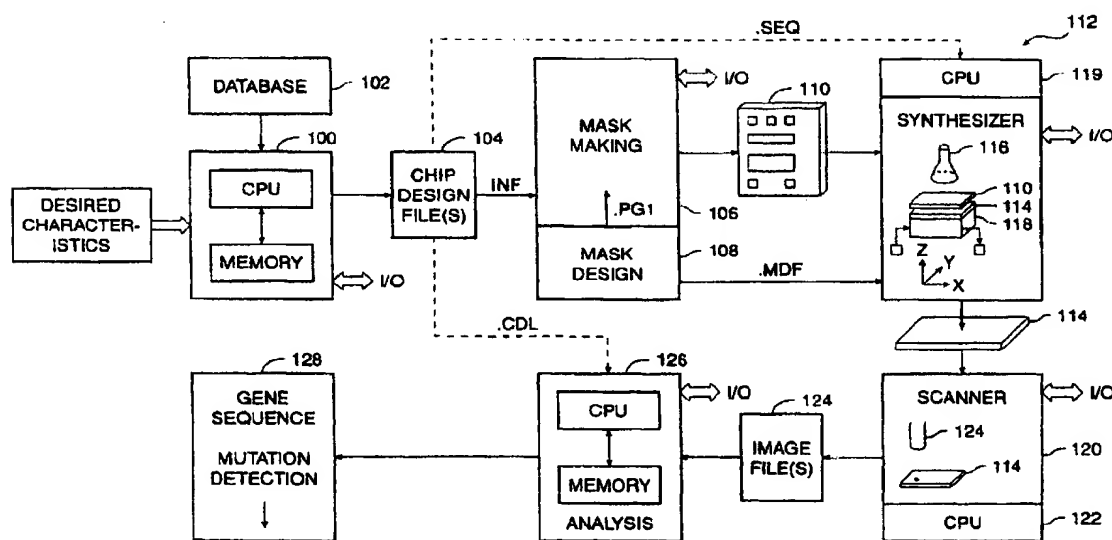


FIG. 3

EP 0 935 210 A3



European Patent
Office

EUROPEAN SEARCH REPORT

Application Number
EP 99 30 0900

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int.Cl.6)
Y,D	WO 97 27317 A (CHEE MARK ;LAI CHAOQIANG (US); LEE DANNY (US); AFFYMETRIX INC (US)) 31 July 1997 (1997-07-31) * the whole document *	1-54	606F19/00
Y	US 4 845 653 A (CONRAD MORGAN P ET AL) 4 July 1989 (1989-07-04) * abstract *	1-54	
A	EP 0 561 241 A (IBM) 22 September 1993 (1993-09-22) * abstract *	1-54	
A	WEINSTEIN JOHN N ET AL: "An information-intensive approach to the molecular pharmacology of cancer" SCIENCE, AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE,, US, vol. 275, no. 5298, 17 January 1997 (1997-01-17), pages 343-349, XP002199806 ISSN: 0036-8075 * abstract; figure 5 *	1-54	
			TECHNICAL FIELDS SEARCHED (Int.Cl.6)
			606F
The present search report has been drawn up for all claims			
Place of search THE HAGUE		Date of completion of the search 25 September 2002	Examiner Filloy García, E
CATEGORY OF CITED DOCUMENTS X: particularly relevant if taken alone Y: particularly relevant if combined with another document of the same category A: technological background O: non-written disclosure P: intermediate document		T: theory or principle underlying the invention E: earlier patent document, but published on, or after the filing date D: document cited in the application L: document cited for other reasons A: member of the same patent family, corresponding document	

EPC FORM 1503 03 02 (PstC01)

**ANNEX TO THE EUROPEAN SEARCH REPORT
ON EUROPEAN PATENT APPLICATION NO.**

EP 99 30 0900

This annex lists the patent family members relating to the patent documents cited in the above-mentioned European search report. The members are as contained in the European Patent Office EDP file on
The European Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

25-09-2002

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
WO 9727317 A	31-07-1997	AU 2253397 A	20-08-1997
		EP 0880598 A1	02-12-1998
		JP 2002515738 T	28-05-2002
		WO 9727317 A1	31-07-1997
		US 6344316 B1	05-02-2002
US 4845653 A	04-07-1989	NONE	
EP 0561241 A	22-09-1993	US 6384847 B1	07-05-2002
		CA 2082917 A1	21-09-1993
		EP 0561241 A2	22-09-1993
		JP 2587894 B2	05-03-1997
		JP 6083975 A	25-03-1994

EPO FORM 10459

For more details about this annex : see Official Journal of the European Patent Office, No. 12/82

THIS PAGE BLANK (ISPTC)